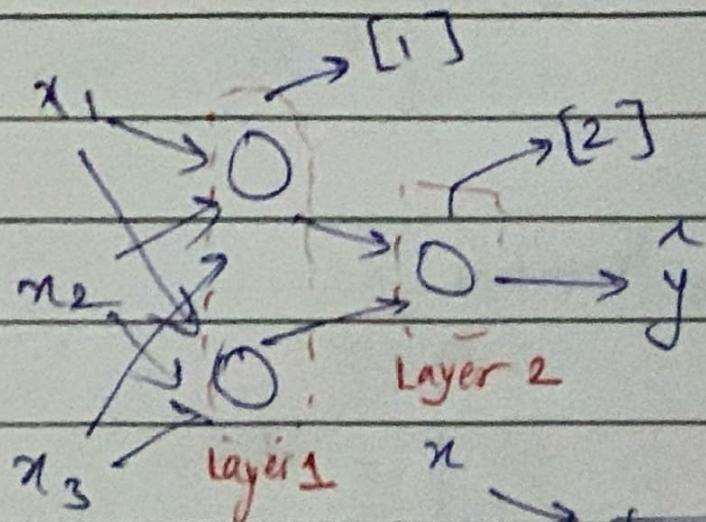


WEEK 3: Shallow Neural Networks

Date: 22/07/24



$$w^{[1]} \rightarrow z^{[1]} = w^{[1]} x + b^{[1]}$$

$$b^{[1]}$$

Further

$$a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$w^{[2]}, b^{[2]}$$

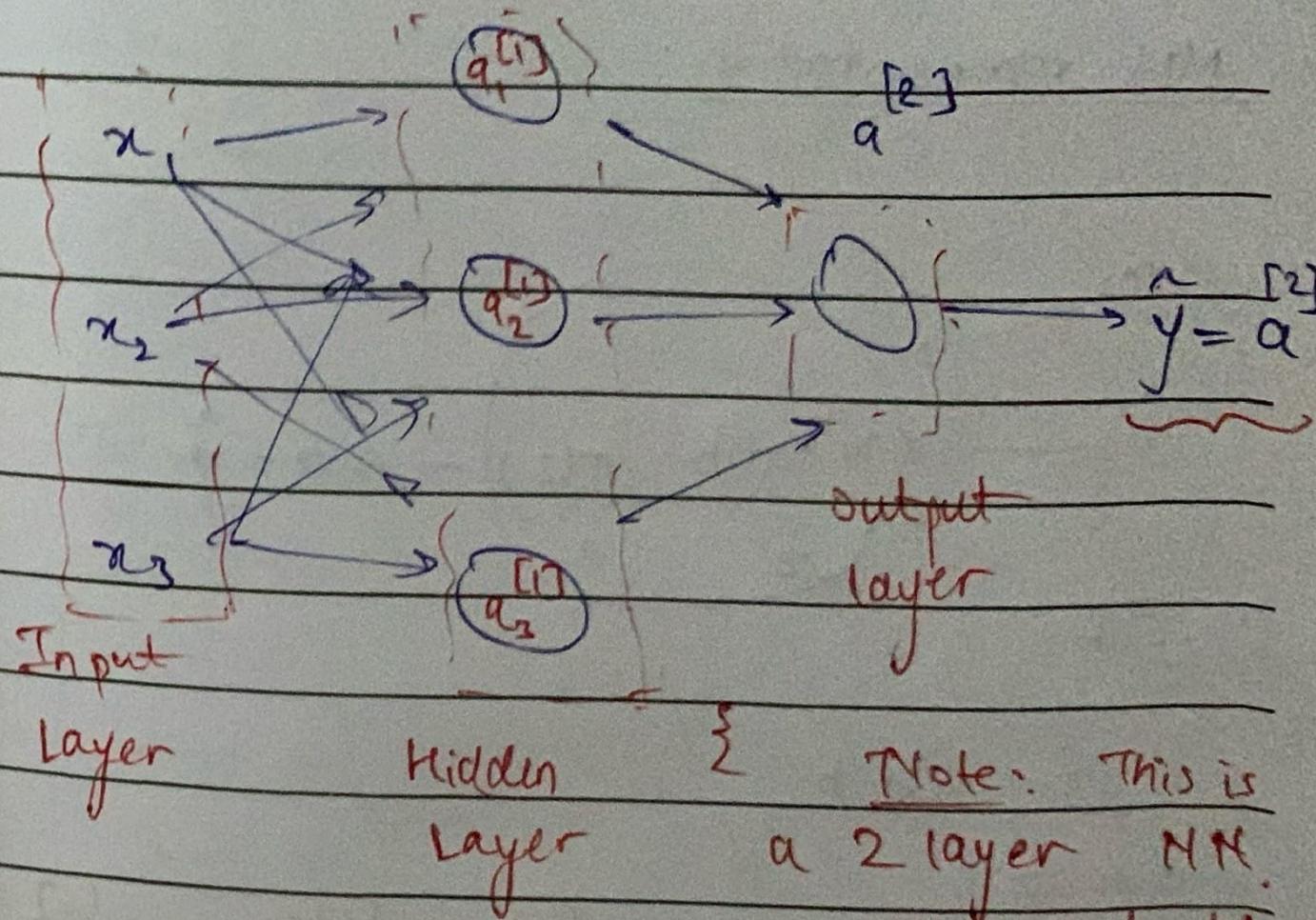
$$a^{[2]} = \sigma(z^{[2]})$$

$$L(a^{[2]}, y)$$

$$a^{[0]} = X$$

$$a^{[1]}$$

Date:

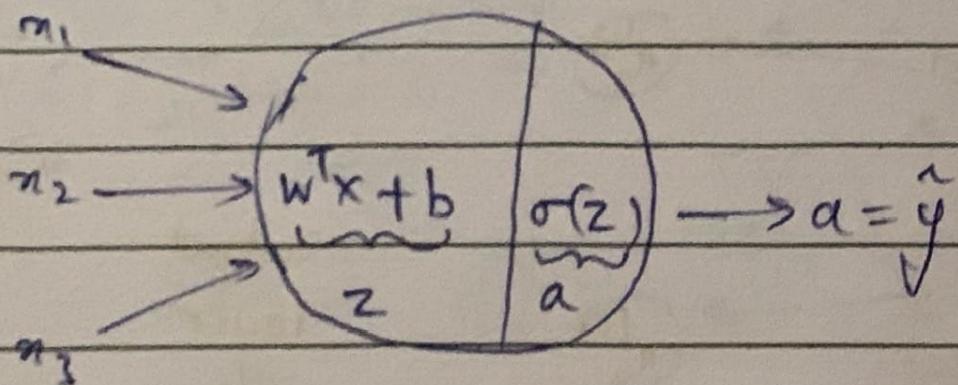


Note: This is
a 2 layer NN.
we don't count the
input layer,

In the hidden layer, you can't see
what's there in the training set, while
you can see it in the input and output
layer.

Date:

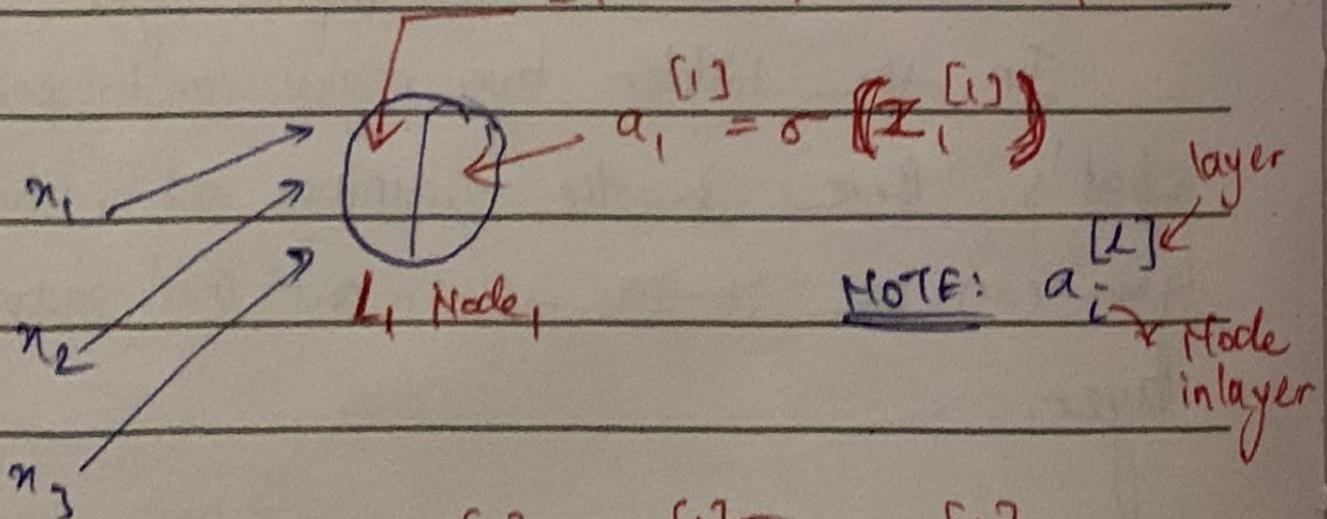
NN representation:



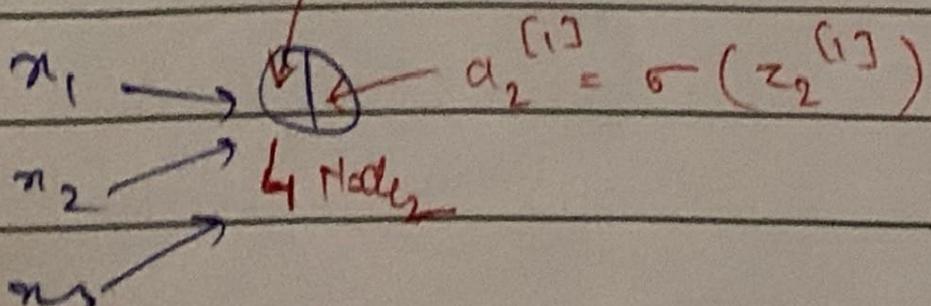
$$z = w^T x + b$$

$$a = \sigma(z)$$

$$z^{(1)} = w_1^{(1)T} x + b_1^{(1)}$$



$$z_2^{(1)} = w_2^{(1)T} x + b_2^{(1)}$$



Date:

Similarly,

For all 4 nodes in layer,
the $z^{(1)}$ is computed.

To vectorise: $W^{(1)}$

$b^{(1)}$

$$Z^{(1)} = \left\{ \begin{array}{c} w_1^{(1)T} \\ w_2^{(1)T} \\ w_3^{(1)T} \\ w_4^{(1)T} \end{array} \right\} \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{array} \right]$$

(4x3)

$$= \left[\begin{array}{c} w_1^{(1)T} x + b_1^{(1)} \\ w_2^{(1)T} x + b_2^{(1)} \\ w_3^{(1)T} x + b_3^{(1)} \\ w_4^{(1)T} x + b_4^{(1)} \end{array} \right]$$

$$= \left[\begin{array}{c} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{array} \right]$$

Date:

$$z^{[1]} = w^{[1]T} x + b^{[1]}$$

$\downarrow \quad \downarrow \quad \downarrow$
 $(4, 1) \quad (4, 3), (3, 1) \quad (1, 1)$

$$a^{[1]} = \sigma(z^{[1]})$$

$\downarrow \quad \downarrow$
 $(4, 1) \quad (4, 1)$

$$z^{[2]} = w^{[2]T} x^{[1]} + b^{[2]} \quad \{w^T = W^{[2]}\}$$

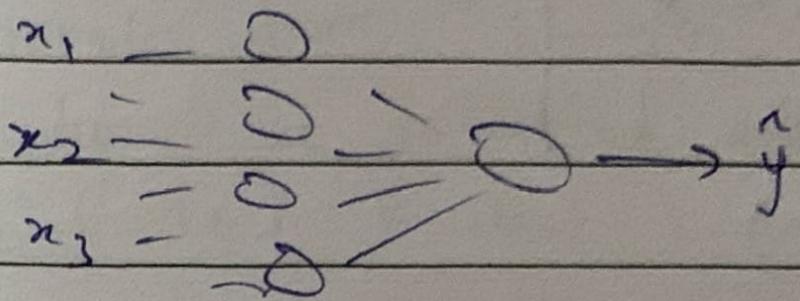
$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $(1, 1) \quad (1, 4) \quad (4, 1) \quad (1, 1)$

$$a^{[2]} = \sigma(z^{[2]})$$

$\downarrow \quad \downarrow$
 $(1, 1) \quad (1, 1)$

Date:

Vectorising across multiple training examples:



$$\begin{aligned} x &\longrightarrow a^{[2]} = \hat{y} \\ x^{(1)} &\longrightarrow a^{[2](1)} = \hat{y}^{(1)} \\ x^{(2)} &\longrightarrow a^{2} = \hat{y}^{(2)} \\ \vdots & \vdots \\ x^{(m)} &\longrightarrow a^{[2](m)} = \hat{y}^{(m)} \end{aligned}$$

$a^{[1](m)}$
layer training
examples

for $i = 1$ to m : {

$$z^{(1)(i)} = w^{(1)} x^{(i)} + b^{(1)}$$

$$a^{(1)(i)} = \sigma(z^{(1)(i)})$$

$$z^{(2)(i)} = w^{(2)} a^{(1)(i)} + b^{(2)}$$

$$a^{(2)(i)} = \sigma(z^{(2)(i)})$$

}

Date:

$$X = \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & | & | \end{bmatrix}$$

(n_n, m)

$$z^{[1]} = w^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

$$z^{[1]} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & \cdots & z^{(1)(m)} \\ | & | & | & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} ; & \{ & \\ a^{(1)(1)} & a^{(1)(2)} & \cdots & a^{(1)(m)} \\ | & | & | & | \end{bmatrix}$$

Date: _____

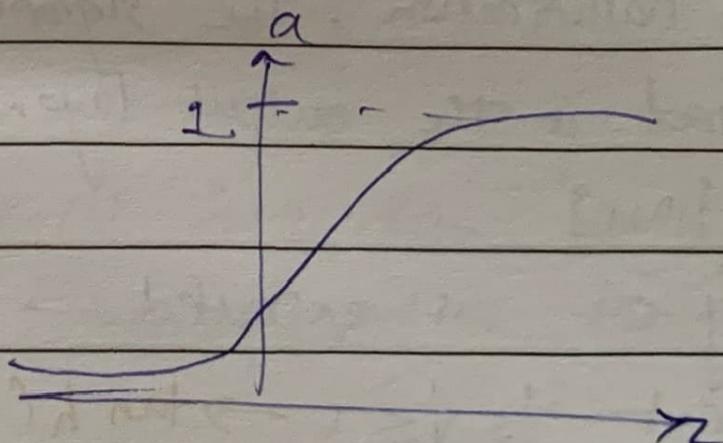
horizontally, we go across different training examples from $1 \rightarrow m$

vertically, we go across different nodes of layer $[l]$ that $1 \rightarrow \text{no. of output nodes}$
 $(\rightarrow \text{no. of hidden units})$

Activation functions :

We used sigmoid function until now, but it is not the best choice.

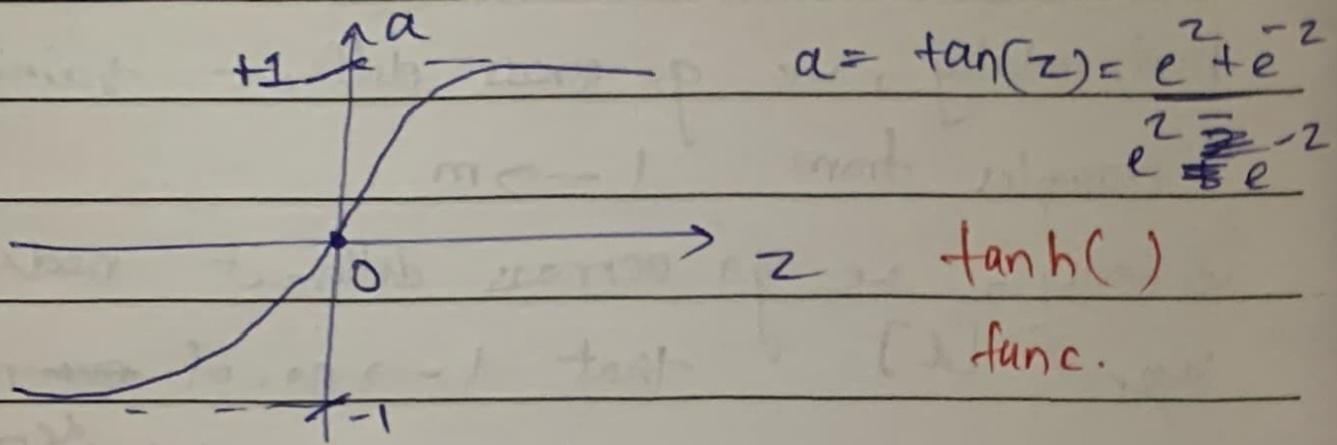
Sigmoid is a type of activation function.



$$a = \frac{1}{1 + e^{-z}} = \sigma(z)$$

Sigmoid func.

Date:



This results better by centering your data by that the mean of your data is close to zero rather than maybe 0.5 , makes learning for next layer a bit easier. $\tanh(x)$ always better than $\sigma(x)$.

Still during binary classification, the sigmoid func. must be used in the output layer

$\therefore \text{inter} \rightarrow y \{0, 1\}$

Sigmoid $\leftarrow 0 \leq \hat{y} \leq 1$ is expected

and not $-1 \leq \hat{y} \leq 1 \rightarrow \tanh()$

Similar problem amongst both is that is ' z ' is very large (or) very small then the slope is almost zero and this slows down gradient descent.

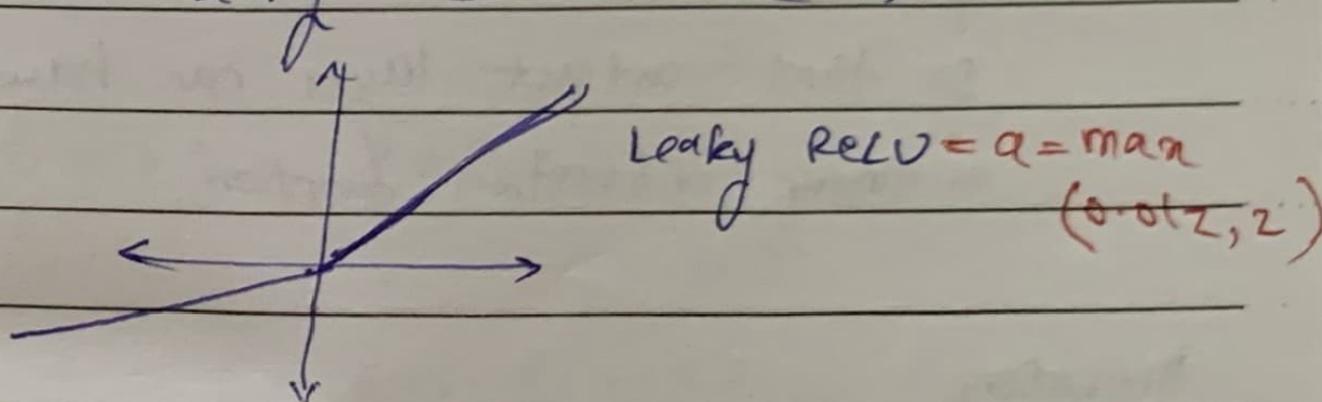
Date:

If binary classification (0/1) \rightarrow sigmoid activation function for output layer

If for hidden layers,

nothing is defined, it is most often ReLU being used (or sometimes $\tanh()$)
problem with ReLU is that for negative values of z , the slope is zero.

L. Leaky ReLU can be used.



ReLU faster than ' \tanh ' and 'sigmoid'
 \therefore learning will be faster.

ReLU / Leaky ReLU $>$ $\tanh()$ $>$ sigmoid

Date:

Activation function :

These are to be non-linear functions, otherwise, there is no point of the hidden layer.

If $y \in \mathbb{R}$

then hidden layers could take up ReLU, tanh so that output layer can have

a linear activation function.

Derivatives

Derivatives :

(1)

$$\frac{d}{dz} g(z)$$

$g = \text{sigmoid funct.}$

$$= \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) = \frac{-1}{(1+e^{-z})^2} \cdot \frac{1}{1+e^{-z}} \cdot e^{-z}$$

Date:

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{1+e^{-z}-1}{(1+e^{-z})^2}$$

$$= \frac{1}{1+e^{-z}} - \frac{1}{(1+e^{-z})^2}$$

$$= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$\Rightarrow \frac{d}{dz} g(z) = \left(g(z) \right) \left(1-g(z) \right)$$

if $z=10$, $g(z) \approx 1$

$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$$

if $z=-10$, $g(z) \approx 0$

$$\frac{d}{dz} g(z) \approx 0 \times (1-0) \approx 0$$

Date:

$$\text{if } z=0 \quad g(z) = y_2$$

$$\frac{d}{dz} g(z) = \frac{1}{2} (1 - y_2) = y_4$$

$$g'(z) = \frac{d}{dz} g(z)$$

$$\text{if } a = g(z)$$

$$\therefore g'(z) = a(1-a)$$

② $\tanh(\cdot)$ activation function

$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$

$$g(z) = \tanh(z) = \frac{e^z + e^{-z}}{e^z - e^{-z}}$$

$$g(z) = 1 - (\tanh(z))^2$$

Date:

$$z = 10 \quad \tanh \approx 1$$

$$g'(z) \approx 0$$

$$z = -10 \quad \tanh \approx -1$$

$$g'(z) \approx 0$$

$$z = 0 \quad \tanh \approx 0$$

$$g'(z) = 1$$

$$a = g(z), \quad g'(z) = \frac{1}{1-a^2}$$

③ ReLU and Leaky ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Date:

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Gradient descent for ANN: $\left[\begin{array}{l} n_2 = n^{[0]}, n^{[1]}, \\ n^{[2]} \end{array} \right] \Rightarrow 1$

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$$(n^{[1]}, n^{[0]}) (n^{[1]}, 1) (n^{[2]}, n^{[1]}) (n^{[2]}, 1)$$

Cost function:

$$\mathcal{J}(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$$

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y)$$

$$a^{[2]}$$

Date:

Gradient descent:

Repeat {

Compute predict ($\hat{y}^{(i)}$, $i=1, 2, \dots, m$)

$$dw^{(1)} = \frac{\partial J}{\partial w^{(1)}}, db^{(1)} = \frac{\partial J}{\partial b^{(1)}}, \dots$$

$$w^{(1)} = w^{(1)} - \alpha dw^{(1)}$$

$$b^{(1)} = b^{(1)} - \alpha db^{(1)}$$

}

}

Date:

formulas for computing derivatives:

→ Forward propagation:

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

→ Back propagation:

$$y = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}]$$

$$dz^{[2]} = A^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]} T$$

$(n^{[1]}, m)$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, axis=1, keepdims=True)$$

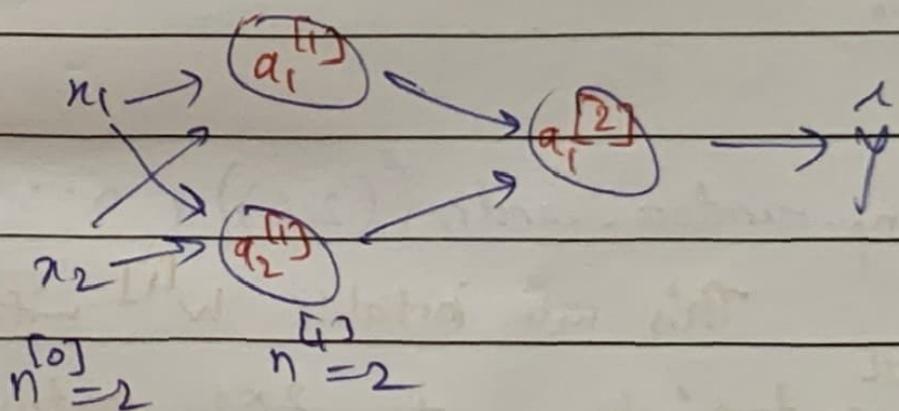
Date:

$$dz^{[1]} = \underbrace{w^{[2]T} dz^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{\text{element-wise product}} (n^{[1]}, m)$$

$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum} (dz^{[1]}, \text{axis}=1, \text{keepdim}=\text{True})$$

Random initialisation:



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Date:

$$a_1^{[1]} = a_2^{[1]} ; \quad dz_1^{[1]} = dz_2^{[1]}$$

$$w^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix} \quad dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

every is same

\therefore symmetric

The NN having hidden units all will compute the same function in no point.

→ Random:

$$w^{[1]} = np.random.randn(2, 2) * 0.01$$

This will initialise $w^{[1]}$ with very small random values of 2×2 matrix

$$b^{[1]} = np.zeros((2, 1))$$

$w^{[2]}$ → randomly

$b^{[2]}$ → to zeros.

Date:

The constant w_{cd} is '0.01' very small

\therefore if activation function used is:

$\tanh(\cdot)$ (or) sigmoid ; then the large values will give almost zero slope and gradient descent will be very ~~slow~~ slow so the model will learn ^{very} slowly.